- 28 -

CLAIMS:

1.    A method of searching for a plurality of query
sequences in a set of target sequence fragments, allowing
5    for mismatches at up to n sequence positions, comprising the
steps of:
     a) logically dividing each query sequence and each
target fragment into at least n+1 segments;
     b) for each query sequence, constructing a first and a
10   second query group by distributing the query sequence
segments there between such that at least n segments are
contained in the second query group;
     c) constructing from each target fragment a first
target group having the same segment distribution as the
15   first query group; and
     d) for each query sequence, comparing said first query
group with each first target group to identify potential
matching target fragments.

20   2.    The method of claim 1 wherein steps (b) to (d) are
repeated using different distributions of segments between
said query groups.

3.    The method of claim 2 wherein steps (b) to (d) are
25   repeated using all distinct distributions of segments
between said query groups.

4.    The method of any preceding claim further comprising,
for each potential matching target fragment identified in
30   step (d), the following steps:

- 29 -

(e) constructing a second target group having the same segment distribution as the second query group of the potentially matching query sequence; and

(f) comparing said second query group with said second
5   target group to identify a match, allowing for mismatches in up to n sequence positions.

5.    The method of claim 4 wherein step (f) is carried out by applying an exclusive OR operation to binary
10   representations of the second query group and second target group.

6.    The method of claim 5 wherein the result of said exclusive OR operation is analysed using a lookup table.

15

7.    The method of any preceding claim further comprising, for each distinct distribution of query sequence segments, constructing a first query table indexed by possible values of the first query groups, the entries in the first query
20   table providing access to each second query group by using as an index the value of a corresponding first query group.

8.    The method of claim 7 further comprising, for each distinct distribution of query sequence segments,
25   constructing a second query table providing access to each second query group, the entries in the first query table providing references to appropriate entries in the second query table.

30   9.    The method of claim 7 or 8 further comprising the step of, for each first target group constructed in step (c),

- 30 -

implementing step (d) by using said first target group to
form an index into said first query table.

10.   The method of claim 7, 8 or 9 wherein, if two distinct

5    distributions of query sequence segments are such that the
first query group of one distribution is the same as the
second query group of the other distribution, said query
tables for both of the distinct distributions are
constructed and/or used concurrently.

10

11.   The method of any preceding claim wherein said target
sequence fragments comprise overlapping fragments of one or
more target sequences.

15   12.   The method of any preceding claim wherein said query
sequences and target sequence fragments comprise biochemical
sequence data.

13.   The method of claim 12 wherein said query sequences and

20   target sequence fragments are binary encoded.

14.   The method of any preceding claim wherein $n$ is at least
two.

25   15.   The method of any preceding claim wherein said query
sequences and target sequence fragments are logically
divided into an even number of segments, and the segments
are distributed in equal numbers between the first and
second query and target groups.

30

16.   Computer apparatus comprising elements adapted to carry
out the method of any preceding claim.

17.   A computer program product carrying computer program code comprising elements adapted to carry out the method of any of claims 1 to 15 when executed on a computer.

5

18.   A computer readable medium comprising computer program code adapted to carry out the method steps of any of claims 1 to 15 when executed on a computer.

10   19.   Computer apparatus comprising a memory, for searching for a plurality of query sequences in a set of target sequence fragments, allowing for mismatches at up to **n** sequence positions comprising:
        a query groups constructor adapted to construct in the
15   memory, for each query sequence, first and second query groups by logically dividing each query sequence into **n+1** segments and distributing the query sequence segments between the first and second query groups in one or more ways such that at least **n** segments are contained in each
20   second query group;
        a target groups constructor adapted to construct in the memory, for each target sequence fragment, one or more first target groups having segment distributions corresponding to the first query groups; and
25        a first group comparator adapted, for each query sequence, to compare said one or more first query groups with corresponding ones of said one or more first target groups and to thereby output a result identifying potentially matching query sequences and target fragments.

30

20.   The apparatus of claim 19, wherein the query groups constructor and target groups constructor are adapted to

- 32 -

construct, for each query sequence and target fragment
respectively, a plurality of groups having all distinct
distributions of segments between said first and second
groups.

5

21. The apparatus of claim 19 or 20 wherein the target
groups constructor is adapted to construct, at least for
each target sequence fragment identified by the first group
comparator as a potentially matching target fragment, one or

10   more second target groups having segment distributions
corresponding to the second query groups, and further
comprising a second group comparator adapted, for each
identified potentially matching query group and target
fragment, to compare the corresponding second query group

15   and second target group to identify a match allowing for
mismatches in up to n sequence positions.

22. The apparatus of any of claims 19 to 21, further
adapted, for each distinct distribution of query sequence

20   segments, to construct in the memory a first query table
indexed by possible values of the first query groups, the
entries in the first query table providing access to each
second query group by using as an index the value of a
corresponding first query group.

25

23. The apparatus of claim 22 further adapted, for each
distinct distribution of query sequence segments, to
construct in the memory a second query table providing
access to each second query group, the entries in the first

30   query table providing references to appropriate entries in
the second query table.

- 33 -

24. A computer program product comprising computer program code for searching for a plurality of query sequences in a set of target sequence fragments, allowing for mismatches at up to $n$ sequence positions comprising:

5        a query groups constructor adapted to construct in the memory, for each query sequence, first and second query groups by logically dividing each query sequence into $n+1$ segments and distributing the query sequence segments between the first and second query groups in one or more

10     ways such that at least $n$ segments are contained in each second query group;

a target groups constructor adapted to construct in the memory, for each target sequence fragment, one or more first target groups having segment distributions corresponding to

15     the first query groups; and

a first group comparator adapted, for each query sequence, to compare said one or more first query groups with corresponding ones of said one or more first target groups and to thereby output a result identifying

20     potentially matching query sequences and target fragments.